



国之重器出版工程


网络强国建设

学术中国 · 大数据

Big Data
Processing Platform

大数据处理平台

宋杰 著

 中国工信出版集团


 人民邮电出版社
POSTS & TELECOM PRESS



大数据处理平台

**Big Data
Processing Platform**

宋杰 著



人民邮电出版社
北京

图书在版编目 (C I P) 数据

大数据处理平台 / 宋杰著. — 北京 : 人民邮电出版社, 2017. 12

(学术中国·大数据)

国之重器出版工程

ISBN 978-7-115-46689-1

I. ①大… II. ①宋… III. ①数据处理 IV.
①TP274

中国版本图书馆CIP数据核字(2017)第261361号

内 容 提 要

本书从数据查询、数据分析和迭代计算平台 3 个方面对大数据处理平台的体系结构、基本原理、主流技术、国内外研究进展和成果进行了全面、深入的阐述,对大数据实时处理平台的架构和核心技术进行了展望。企业技术人员可参考本书选择合适的技术构建大数据处理平台或对现有平台进行优化;高校及科研院所的科研人员可参考本书了解大数据管理的基本原理和现有研究成果;高校学生可通过学习本书全面了解大数据处理平台。同时,本书也适用于对大数据技术有浓厚兴趣的读者。

◆ 著 宋 杰

责任编辑 吴娜达

责任印制 杨林杰

◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路 11 号

邮编 100164 电子邮件 315@ptpress.com.cn

网址 <http://www.ptpress.com.cn>

固安县铭成印刷有限公司印刷

◆ 开本: 720×1000 1/16

印张: 14.75

2017 年 12 月第 1 版

字数: 220 千字

2017 年 12 月河北第 1 次印刷

定价: 89.00 元

读者服务热线: (010)81055488 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京东工商广登字 20170147 号

《国之重器出版工程》 编辑委员会

编辑委员会主任：苗 圩

编辑委员会副主任：刘利华 辛国斌

编辑委员会委员：

| | | | | |
|-----|-----|-----|-----|-----|
| 冯长辉 | 梁志峰 | 高东升 | 姜子琨 | 许科敏 |
| 陈 因 | 郑立新 | 马向晖 | 高云虎 | 金 鑫 |
| 李 巍 | 李 东 | 高延敏 | 何 琼 | 刁石京 |
| 谢少锋 | 闻 库 | 韩 夏 | 赵志国 | 谢远生 |
| 赵永红 | 韩占武 | 刘 多 | 尹丽波 | 赵 波 |
| 卢 山 | 徐惠彬 | 赵长禄 | 周 玉 | 姚 郁 |
| 张 炜 | 聂 宏 | 付梦印 | 季仲华 | |




专家委员会委员（按姓氏笔画排列）：

- 于 全 中国工程院院士
- 王少萍 “长江学者奖励计划”特聘教授
- 王建民 清华大学软件学院院长
- 王哲荣 中国工程院院士
- 王 越 中国科学院院士、中国工程院院士
- 尤肖虎 “长江学者奖励计划”特聘教授
- 邓宗全 中国工程院院士
- 甘晓华 中国工程院院士
- 叶培建 中国科学院院士
- 朱英富 中国工程院院士
- 朵英贤 中国工程院院士
- 邬贺铨 中国工程院院士
- 刘大响 中国工程院院士
- 刘怡昕 中国工程院院士
- 刘韵洁 中国工程院院士
- 孙逢春 中国工程院院士
- 苏彦庆 “长江学者奖励计划”特聘教授



- 苏哲子 中国工程院院士
- 李伯虎 中国工程院院士
- 李应红 中国科学院院士
- 李新亚 国家制造强国建设战略咨询委员会委员、
中国机械工业联合会副会长
- 杨德森 中国工程院院士
- 张宏科 北京交通大学下一代互联网互联设备国家
工程实验室主任
- 陆建勋 中国工程院院士
- 陆燕荪 国家制造强国建设战略咨询委员会委员、原
机械工业部副部长
- 陈一坚 中国工程院院士
- 陈懋章 中国工程院院士
- 金东寒 中国工程院院士
- 周立伟 中国工程院院士
- 郑纬民 中国计算机学会原理事长
- 郑建华 中国科学院院士

- 
- 屈贤明 国家制造强国建设战略咨询委员会委员、工业和信息化部智能制造专家咨询委员会副主任
- 项昌乐 “长江学者奖励计划”特聘教授，中国科协书记处书记，北京理工大学党委副书记、副校长
- 柳百成 中国工程院院士
- 闻雪友 中国工程院院士
- 徐德民 中国工程院院士
- 唐长红 中国工程院院士
- 黄卫东 “长江学者奖励计划”特聘教授
- 黄先祥 中国工程院院士
- 黄 维 中国科学院院士、西北工业大学常务副校长
- 董景辰 工业和信息化部智能制造专家咨询委员会委员
- 焦宗夏 “长江学者奖励计划”特聘教授

《学术中国·大数据》丛书 编辑委员会

编辑委员会顾问：

邬贺铨 李国杰 李德毅 方滨兴

编辑委员会主任：郑纬民

编辑委员会委员（按姓氏笔画排列）：

王建民 杜跃进 李国庆 李 涛 宋 杰

张广艳 陈 卫 陈世敏 魏哲巍

策 划：《大数据》杂志



丛书总序

大数据、人工智能、云计算、物联网、移动互联网和产业互联网等成为新一代信息技术的特征，其中大数据与上述技术和应用都有密切关系。大数据来自于移动互联网、产业互联网和物联网等，其存储需要云计算，其挖掘依靠人工智能，而人工智能也有赖于大数据的支撑，大数据是产业互联网的重要基础。大数据不仅可以用于社会的精细化管理，更好地服务民生，大数据产业也将形成信息产业新的分支，其间接的产业影响将更大。可以说，大数据是数字经济的重要支柱。

很多国家都将大数据作为新时期的国家发展战略。2015年，国务院印发大数据发展的首个权威性、系统性文件《促进大数据发展行动纲要》，2016年国家发展和改革委员会批复了13个大数据领域的国家工程实验室，我国一些省市也纷纷制定大数据发展战略与规划。当前，我国在大数据共享开放、大数据资源开发、大数据技术研发、大数据挖掘应用、大数据产业培育、大数据安全管理、大数据人才培养和大数据法规研究等方面全面部署，为我国实现供给侧结构性改革，促进产业升级和转型，提升国家竞争力，争取在国际领域的话语权和实现跨越式发展起到了不可或缺的作用。

然而，我国的大数据发展也面临一些亟待解决的问题，例如基础研究薄弱、创新能力不强、产业链条缺口、数据资源封闭、法律法规滞后、数据安全不力、数据人才短缺和数据设施布局不合理及利用率不高等。为了使我国的大数据应用与产业可持续健康发展，需要多管齐下，其中普及大数据科学是重要的一环。为此，《学术中国·大数据》丛书编委会组织多个大数据领域优秀的研究团队的专家，基于国家



“973”计划、“863”计划、国家自然科学基金、国家重点研究计划等科研项目的创新研究成果和国内外大数据应用的成功实践，编写了这套丛书，内容涵盖大数据存储、数据管理、数据挖掘、分析平台、优化算法等核心技术领域。

本丛书的出版对传播大数据科学知识、推动大数据的学术探讨、鼓励大数据领域的产学研用协同创新、促进大数据标准化研究、加快大数据核心技术研发、培训大数据技术人才、引导大数据应用与产业化发展以及完善大数据有关的制度建设，都将起到积极作用。

2017年12月



前言

如何从海量数据中有效获取信息，以进行分析和决策是大数据的核心问题，也是 21 世纪各行各业均面临的重要问题。解决这一核心问题需要大数据处理平台的支持。大数据处理平台是一种“计算平台”，计算平台泛指支持算法执行的硬件系统、操作系统和运行库，那么大数据处理平台则泛指可以支持大数据处理算法执行的平台。大数据处理平台采用集群作为硬件系统，分布式计算框架作为中间件系统。以 Hadoop HDFS 为代表的分布式文件系统和以 MapReduce 为代表的分布式并行编程模型在学术界和产业界最为流行，并以此引出完善的 Hadoop 生态圈；另一个则是围绕并行框架 Spark 的生态系统，如 Spark Streaming 和 Shark。以这些开源技术为支撑的大数据处理平台广泛地应用于社交网络、科学数据分析、传感数据处理、医疗和电子商务平台中。

典型的大数据应用可以分为 OLTP、OLAP 和图计算 3 类，因此，从数据处理平台角度，需要提供数据查询、统计分析和迭代计算支持。本书围绕这 3 个典型数据处理方式，首先介绍大数据处理平台的体系结构，并简述体系结构每部分的主流技术；随后重点介绍大数据处理平台实现数据查询、统计分析和迭代计算的基本原理、研究进展；每部分还包括项目组近年来的研究成果。此外，本书还展望了大数据实时处理平台的架构和核心技术。本书介绍的理论和技術均集中于中间件层，以“学术研究”和“系统实现”相结合的角度论述，使得读者能够更加深入地理解大数据处理平台的核心技术和学术前沿，帮助读者更加有效地构建处理平台，或对已有的大数据处理平台进行改进，开展大数据存储和管理领域的相关研究。



本书是“学术中国·大数据”系列丛书之一。书中研究成果为国家自然科学基金重点项目“大数据高效能存储与管理方法研究(No.61433008)”的部分建设成果。笔者在数据管理领域已有十多年的研究经验,结合自身的研究经验,从“学术研究”和“系统实现”相结合的角度,对平台进行全面的介绍。书中既有原理,又有学术前沿综述,但不包含使用方法、编程技术、构建步骤等类似工具书的内容。对于大数据相关领域高校师生、研究人员以及大数据处理平台的设计师和架构师有一定的借鉴性。

笔 者

2017年4月于沈阳南湖



目 录

| | |
|------------------------|-----|
| 第 1 章 体系结构 | 001 |
| 1.1 集群系统 | 002 |
| 1.1.1 Hadoop YARN | 002 |
| 1.1.2 Apache Mesos | 003 |
| 1.1.3 Apache ZooKeeper | 004 |
| 1.2 文件系统 | 005 |
| 1.2.1 Google 分布式文件系统 | 006 |
| 1.2.2 Hadoop 分布式文件系统 | 008 |
| 1.2.3 其他分布式文件系统 | 009 |
| 1.3 NoSQL 和 NewSQL | 012 |
| 1.3.1 NoSQL 数据库系统 | 012 |
| 1.3.2 NewSQL 数据库系统 | 014 |
| 1.4 计算模型 | 016 |
| 1.4.1 MapReduce 编程模型 | 016 |
| 1.4.2 Spark 并行计算框架 | 025 |
| 参考文献 | 026 |
| 第 2 章 查询平台 | 031 |
| 2.1 基本原理 | 032 |
| 2.1.1 系统简介 | 033 |
| 2.1.2 架构组织 | 034 |
| 2.2 现有研究 | 037 |
| 2.2.1 大数据精确查询系统 | 037 |
| 2.2.2 大数据近似查询系统 | 040 |
| 2.2.3 大数据多维查询系统 | 040 |



| | |
|------------------------|------------|
| 2.3 近期成果 | 043 |
| 2.3.1 Haery | 043 |
| 2.3.2 Probery | 056 |
| 参考文献 | 075 |
| 第3章 分析平台 | 081 |
| 3.1 基本原理 | 082 |
| 3.1.1 OLAP 技术 | 082 |
| 3.1.2 系统架构 | 084 |
| 3.2 现有研究 | 086 |
| 3.2.1 传统 OLAP 优化方法 | 086 |
| 3.2.2 OLAP 存储计算优化 | 090 |
| 3.2.3 大数据 OLAP 引擎 | 097 |
| 3.3 近期成果 | 098 |
| 3.3.1 DOLAP | 099 |
| 3.3.2 MapReduce OLAP | 109 |
| 3.3.3 HaoLap | 119 |
| 参考文献 | 121 |
| 第4章 迭代计算平台 | 127 |
| 4.1 基本原理 | 128 |
| 4.2 现有研究 | 129 |
| 4.2.1 MapReduce 迭代计算框架 | 130 |
| 4.2.2 其他迭代计算框架 | 132 |
| 4.2.3 增量迭代计算 | 136 |
| 4.2.4 迭代算法优化 | 137 |
| 4.3 近期成果 | 139 |
| 4.3.1 增量迭代计算模型 | 139 |
| 4.3.2 归并迭代计算 | 157 |
| 4.3.3 迭代初始点选择 | 159 |
| 参考文献 | 172 |
| 第5章 实时处理平台 | 175 |
| 5.1 基本原理 | 176 |



| | |
|-----------------------|-----|
| 5.2 现有研究 | 178 |
| 5.2.1 Lambda 架构 | 179 |
| 5.2.2 队列 | 181 |
| 5.2.3 流处理 | 183 |
| 5.2.4 数据流处理框架 | 189 |
| 5.3 近期成果 | 200 |
| 5.3.1 实时数据迁移模型 | 201 |
| 5.3.2 数据源层的优化方法 | 207 |
| 5.3.3 迁移系统设计 | 216 |
| 参考文献 | 218 |
| 后记 | 219 |



第1章 体系结构

大数据处理平台是一种“计算平台”，计算平台泛指支持算法执行的硬件系统、操作系统和运行库，大数据处理平台泛指可以支持大数据处理算法执行的平台。大数据处理平台采用集群作为硬件系统，分布式计算框架作为中间件系统。本章主要介绍大数据处理平台的体系结构，并简述体系结构每部分的主流技术。

| 1.1 集群系统 |

集群系统通常是指构建在计算机集群之上的系统，其将一定数量的计算机连接起来构成分布式系统，作为单独的、统一的计算资源，为上层应用提供计算服务。集群系统对同构或异构的计算资源进行统一管理、调度以及实现分布式协调（Distributed Coordination）。本节主要介绍当前主流的任务管理与调度机制 Hadoop YARN、Apache Mesos 以及分布式协调机制 Apache ZooKeeper。

1.1.1 Hadoop YARN

由于 MapReduce 的 JobTracker/TaskTracker 机制在可扩展性、内存消耗、线程模型、可靠性等方面存在缺陷，且其维护成本很高，因此，Apache 提出了 Hadoop YARN 这种新型的 Hadoop MapReduce 框架。YARN 架构如图 1-1 所示。YARN^[1] 提供强大的资源管理功能，将 JobTracker 分解为两个独立的服务：全局的资源管理器（Resource Manager, RM）和应用程序特有的应用程序管理器（Application Master, AM）。YARN 采用 Master-Slave（主-从）架构，主节点部署 RM，负责整个系统的资源管理和分配。每个从节点上均部署节点资源管理器（Node Manager, NM），RM 对各个 NM 上的资源进行统一管理和调度。